

# Agentic Data Architecture (ADA): Eliminating the API Layer for Hallucination- Free, Sub-100ms Enterprise AI Agents

Nirmal Nambiar

SuperManager AGI Research Lab

## ABSTRACT

Every major framework for deploying AI agents in production shares one unexamined assumption: **data must be accessed through a network boundary**. Whether via the Model Context Protocol (MCP), CLI wrappers, or direct REST APIs, the agent must leave its execution context, traverse a network, authenticate, and re-enter reasoning for every data lookup. This boundary is the structural root cause of the three most damaging enterprise AI failure modes: hallucination (15-25% error rate, \$250M annual industry loss), prohibitive latency (200-500ms per API round-trip), and single-agent throughput ceilings that prevent elastic scaling.

We present **Agentic Data Architecture (ADA)**, a framework that eliminates this assumption through three mutually reinforcing mechanisms: (i) **direct polyglot database connectivity**, native async connections to PostgreSQL, MongoDB, and Redis bypassing all HTTP intermediaries; (ii) **per-subtask RAG grounding**, Sentence-BERT embeddings and HNSW vector retrieval invoked inside each specialist agent's reasoning loop; and (iii) **hierarchical multi-agent orchestration**, a Controller Agent that decomposes queries into a DAG of subtasks dispatched to stateless Specialist Agents via a work-stealing scheduler.

Evaluation across 10,000 enterprise queries demonstrates: **70% hallucination reduction** (22.4% to 4.2%), **81.4% latency reduction** (350ms to 65ms), **6.7x throughput gain** at ten parallel agents, and **92.5% retrieval accuracy**. All results on consumer-grade hardware (Intel i7, 16GB RAM), no GPU, no cloud API. Data lives where agents reason.

**Keywords:** Agentic AI, Direct Database Access, Hallucination Mitigation, Multi-Agent Orchestration, Retrieval-Augmented Generation, Low-Latency Inference, Enterprise AI, Work-Stealing Scheduler, Beehive Architecture, Polyglot Persistence

## 1. Introduction

The deployment of large language models (LLMs) as autonomous agents has moved from research prototypes to enterprise production at a pace that outstrips architectural best practices. By 2025, AI agents were automating project management, financial analysis, HR operations, and customer support at scale. The architectural question is no longer whether to deploy agents, but how to make them reliable, fast, and cost-effective enough to trust with mission-critical enterprise data.

The debate over agent-data integration crystallised publicly in early 2026. Garry Tan (Y Combinator) argued that the Model Context Protocol "eats too much context window" and advocated for lightweight CLI wrappers [1]. The CTO of Perplexity AI announced an internal pivot away from MCP toward direct APIs and CLIs [2]. Security researchers immediately identified the fatal flaw: CLI-to-API integrations cannot propagate agent identity across multi-hop call chains, breaking policy enforcement and enabling impersonation attacks [3]. Their conclusion: implement security correctly and you have reinvented MCP.

**Both sides are fighting over the wrong layer.** The dispute is about tool-calling protocols. The real problem is the **network boundary** that every protocol presupposes. Whether the agent calls a tool via MCP, a CLI, or a REST endpoint, it must leave its execution context, traverse a network, authenticate, serialise/deserialise payloads, and re-enter reasoning. Each hop costs 200-500ms. Each hop is a hallucination vector. Each hop is a rate-limit risk.

ADA eliminates the hop entirely. By connecting agents natively to databases within the same deployment boundary, with RAG grounding scoped per subtask and throughput scaled via hierarchical parallelism, ADA turns the network boundary from a design constraint into a design choice. For enterprise deployments where data ownership is internal, the right choice is: no boundary.

## 1.1 Problem Formalisation

Let  $Q$  be a set of enterprise knowledge-intensive queries. For each query  $q$ , an agent system  $S$  must produce a response satisfying three constraints simultaneously:

- **Accuracy:**  $\Pr[\text{hallucination}(r)] < 5\%$ , the enterprise acceptability threshold.
- **Latency:**  $E[\text{latency}(r)] < 100\text{ms}$  for real-time enterprise workloads.
- **Throughput:**  $\text{Throughput}(S)$  scales near-linearly with hardware addition.

We demonstrate that no existing approach satisfies all three simultaneously (Section 2), and that ADA does (Section 5).

## 1.2 Principal Contributions

1. **ADA Framework:** First unified architecture satisfying accuracy, latency, and throughput constraints simultaneously by eliminating the API abstraction layer via native polyglot database connectivity.
2. **Per-Subtask RAG:** Novel invocation of Sentence-BERT and HNSW retrieval inside each specialist agent's reasoning loop, enabling task-scoped grounding that reduces hallucination to 4.2%.
3. **Work-Stealing Scheduler:** Hierarchical controller-specialist topology with dynamic task reallocation achieving 6.7x speedup at  $m=10$  agents.
4. **Security-Compatible Perimeter Access:** Formal argument that ADA's direct DB model satisfies the same security invariants MCP seeks to enforce, without protocol overhead.
5. **Reproducible Benchmark:** 10,000-query evaluation across five enterprise domains on consumer hardware; all configurations and hyperparameters reported.

## 2. The MCP / CLI / API Trilemma

We characterise the three dominant paradigms for agent-data integration and demonstrate that each fails on at least one critical constraint. We term this the MCP/CLI/API Trilemma.

### 2.1 MCP (Model Context Protocol)

MCP provides structured tool-calling with OAuth-based authentication, user consent screens, and agent identity chaining via act claims in JWT tokens [4]. Its security model is formally sound for multi-agent call chains. However, MCP's tool schemas consume thousands of context-window tokens per session, adding protocol-level latency on top of existing API latency, and requiring manual toggling that creates fragile production pipelines [1]. It violates the latency constraint.

### 2.2 CLI Wrappers

Lightweight (~100 LOC), deterministic, zero protocol overhead. But OAuth tokens identify the CLI client, not the invoking agent. In agent-to-agent call chains, the API receives a token with no agent identity field. API providers cannot enforce per-agent policies. Impersonation is structurally possible. This violates the security constraint [3].

### 2.3 Direct APIs (Code Mode)

Deterministic, reviewable. Avoids MCP's context bloat. But satisfying enterprise security requirements: dynamic client registration, OAuth consent scoping, sensitive-action approval, and agent-experience API design produces a system architecturally equivalent to MCP [3]. Network RTT remains. The latency constraint is violated.

### 2.4 The Shared Root Cause

**Proposition 1 (Network Boundary Inevitability):** Any agent-data integration paradigm that locates data in an external service cannot simultaneously satisfy hallucination rate < 5%, average latency < 100ms, and linear throughput scaling, due to irreducible network RTT ( $\geq 200\text{ms}$ ), serialisation overhead, and rate-limit exposure.

ADA refutes this proposition by eliminating its precondition: data is not external.

**Table 1: The MCP / CLI / API Trilemma**

Approach	Latency	Security	Hallucination	Scale	No Boundary
MCP	High	Strong	None	Medium	No
CLI Wrappers	Medium	Weak	None	Low	No
Direct API	Medium	Medium	None	Medium	No
<b>ADA (Ours)</b>	Low	Strong	70% reduction	High	Yes

All three existing paradigms fail on at least one critical constraint. ADA satisfies all five dimensions.

## 3. ADA Framework

### 3.1 Design Axioms

**Definition 1 (Data Locality):** Agent data-access latency must be bounded by local I/O, not network RTT.

**Definition 2 (Grounded Generation):** For every factual claim in an agent response, there exists a retrieved database record such that the claim is derivable from that record under the agent's generation model.

**Definition 3 (Elastic Concurrency):** System throughput scales as  $\geq 0.6 * m * \text{throughput}(1)$  for  $m$  deployed specialist agents, subject only to hardware resource limits.

### 3.2 Four-Layer Reference Architecture

**L1: Interface.** Accepts queries via REST, gRPC, or CLI. Performs intent classification and routes to the Controller Agent.

**L2: Orchestration.** The singleton Controller Agent decomposes queries into a DAG of subtasks, topologically sorts dependencies, and dispatches ready subtasks to Specialist Agents via a work-stealing scheduler that minimises idle time without global coordination.

**L3: Retrieval.** The RAG Engine performs per-subtask Sentence-BERT embedding and HNSW approximate nearest-neighbour search. Context is scoped to each subtask, preventing context-window bloat and improving retrieval precision.

**L4: Data.** Native async connectors to PostgreSQL (asyncpg), MongoDB (motor), and Redis (aioredis) with connection pooling, zero HTTP intermediary, and a Redis LRU response cache (TTL = 300s).

### 3.3 Agent Taxonomy

**Controller Agent (CA):** Singleton. Maintains work queue  $W$  and decomposes query  $Q$  into subtask DAG  $T$  via intent classification. Work-stealing assignment minimises idle time.

**Specialist Agents (SA1...SAn):** Stateless, horizontally scalable. Each SA: (i) receives subtask  $T_i$ ; (ii) invokes the RAG Engine for task-scoped context  $C_i$ ; (iii) constructs and executes DB query to obtain  $D_i$ ; (iv) generates response  $R_i$  grounded in  $C_i$  and  $D_i$ .

**Aggregator:** Synthesises results via evidence-majority voting. Each claim traced to its retrieved chunk; majority vote prevents any single hallucinating agent from corrupting the final output.

### 3.4 Security Analysis

The security concern raised against CLI/API approaches centres on agent identity propagation across multi-hop call chains. ADA resolves this structurally rather than protocolically:

6. **No external boundary:** Direct DB connectivity operates inside the enterprise security perimeter (VPN, IAM, firewall). There is no OAuth token chain because there is no external service call.

7. **Database-native access control:** PostgreSQL row-level security enforces per-user data isolation; MongoDB collection-level ACLs scope access by role, both more precise than OAuth token scoping.

8. **Hybrid deployment:** For cross-enterprise agent-to-agent calls, ADA complements MCP at the inter-enterprise boundary while retaining direct DB access internally, combining the security of MCP with the performance of ADA.

## 4. Technical Implementation

### 4.1 RAG Pipeline

#### 4.1.1 Document Chunking

Input documents are segmented into windows of  $w = 512$  tokens with sliding overlap  $\delta = 50$  tokens. Chunk count  $c = \text{ceil}((L - \delta) / (w - \delta))$  for document of  $L$  tokens. Overlap preserves cross-boundary semantic continuity, critical for multi-hop enterprise reasoning.

#### 4.1.2 Embedding

We employ all-mpnet-base-v2, a Sentence-BERT variant producing  $d = 768$ -dimensional dense vectors, L2-normalised before indexing. This model achieves state-of-the-art performance on semantic textual similarity benchmarks while remaining deployable on CPU.

#### 4.1.3 Vector Indexing

Embeddings are indexed via Hierarchical Navigable Small World (HNSW) graphs using FAISS 1.7.4 with  $\text{efsearch} = 200$ , providing  $O(\log n)$  approximate nearest-neighbour search at recall  $> 0.98$ .

#### 4.1.4 Retrieval and Grounding Filter

Given subtask query  $q_i$ , top- $k = 5$  chunks are retrieved by cosine similarity:  $\text{sim}(q_i, c_j) = (q_i \cdot c_j) / (\|q_i\| \|c_j\|)$ . Chunks with similarity  $< \tau = 0.70$  are discarded, a hard grounding filter that prevents the LLM from generating beyond retrieved evidence.

**Table 2: RAG Pipeline Configuration**

Parameter	Value
Chunk window ( $w$ )	512 tokens
Chunk overlap ( $\delta$ )	50 tokens
Embedding model	all-mpnet-base-v2
Vector dimension	768
Normalisation	L2
Index algorithm	HNSW
efsearch	200
Recall@10	$> 0.98$
Top-k	5
Similarity threshold ( $\tau$ )	0.70

Table 2: RAG Pipeline Configuration

## 4.2 Multi-Agent Orchestration Algorithm

The hierarchical orchestration procedure is as follows:

INPUT: Query Q; agent pool  $A = \{A_1, \dots, A_m\}$

1. CA classifies intent of Q
  2.  $T \leftarrow \text{Decompose}(Q)$  [DAG of subtasks]
  3.  $W \leftarrow \text{TopologicalSort}(T)$
  4. WHILE W is not empty:
    - $T_{\text{ready}} \leftarrow \{T \text{ in } W : \text{deps}(T) = \text{empty}\}$
    - FOR ALL  $T_i$  in  $T_{\text{ready}}$  IN PARALLEL:
      - $A_i \leftarrow \text{WorkStealAssign}(T_i, A)$
      - $C_i \leftarrow \text{RAGEngine.Retrieve}(T_i)$  [per-subtask]
      - $q_i \leftarrow \text{BuildDBQuery}(T_i, C_i)$
      - $D_i \leftarrow \text{DBConnector.Execute}(q_i)$
      - $R_i \leftarrow \text{LLM.Generate}(T_i, C_i, D_i)$
    - Remove  $T_i$  from W; update dependency graph
  5.  $R \leftarrow \text{EvidenceMajorityAggregate}(\{R_i\})$
- OUTPUT: Grounded response R

## 4.3 Database Connectivity Layer

**Table 3: Database Connector Performance and Access Controls**

Database	Driver	Latency	Throughput	Access Control
PostgreSQL	asyncpg	15ms	800 q/s	Row-level security
MongoDB	motor	8ms	1200 q/s	Collection-level ACL
Redis	aioredis	2ms	5000 q/s	AUTH + TLS + TTL

Table 3: Database Connector Performance. All connectors use parameterised queries. Redis serves as LRU response cache (TTL = 300s).

## 5. Evaluation

### 5.1 Experimental Setup

**Hardware:** Intel Core i7-12700K (12 cores, 3.6GHz), 16GB DDR5-4800, 1TB NVMe SSD. No GPU. No cloud API. All database instances co-located on the same host.

**Software:** Python 3.11, PyTorch 2.2 (CPU), HuggingFace Transformers 4.38, FAISS 1.7.4, asyncpg 0.29, motor 3.3, aioredis 2.0, LangChain 0.1, Llama-2-7B (4-bit GGUF quantisation, llama.cpp backend).

**Baselines:** (1) Vanilla LLM: Llama-2-7B, no retrieval, no grounding. (2) API-RAG: Same LLM, RAG over Bing Search API v7. (3) Single-Agent ADA: Full ADA pipeline,  $m=1$ . (4) ADA (Proposed): Full framework,  $m$  in  $\{1, 3, 5, 10\}$ .

**Benchmark:** 10,000 queries across five enterprise domains: customer support (2k), financial analytics (2k), HR intelligence (2k), product knowledge (2k), and operational data (2k). 30% require multi-hop reasoning across 3-7 database lookups. Ground-truth answers curated by domain experts; hallucination labelled by two independent annotators ( $\kappa = 0.87$ ).

**Metrics:** Hallucination Rate (HR, %), Average Latency (ms), Throughput (tasks/s), Retrieval Accuracy (RA, %), P95 Latency (ms). Statistical test: paired t-test,  $p < 0.001$ ,  $n = 10,000$ ; effect sizes reported as Cohen's  $d$ .

## 5.2 System-Level Results

**Table 4: End-to-End System Comparison (10,000 Queries)**

System	HR (%)	Avg Latency (ms)	RA (%)	P95 (ms)
Vanilla LLM	22.4	310	61.2	520
API-RAG	14.1	350	76.8	610
Single-Agent ADA	5.8	78	88.3	140
<b>ADA (10 agents)</b>	<b>4.2</b>	<b>65</b>	<b>92.5</b>	<b>112</b>

Table 4: All ADA vs. API-RAG improvements:  $p < 0.001$ , Cohen's  $d > 1.2$  (large effect). HR = Hallucination Rate; RA = Retrieval Accuracy.

## 5.3 Latency Decomposition

**Table 5: Latency Decomposition by Component (ms)**

System	Embed (te)	Retrieval (tr)	Database (td)	Generation (tg)	Total
API-RAG	12	285	n/a	53	350
ADA (1 agent)	12	8	15	43	78
<b>ADA (10 agents)</b>	<b>12</b>	<b>8</b>	<b>8</b>	<b>37</b>	<b>65</b>

Table 5: Eliminating network-bound API retrieval (285ms) is the dominant saving, confirming Data Locality (Definition 1) as the primary latency driver.

## 5.4 Multi-Agent Scaling

**Table 6: Multi-Agent Scaling Results (500 Concurrent Tasks)**

Agents (m)	Time (s)	Speedup	Efficiency
1	500	1.0x	100%

Agents (m)	Time (s)	Speedup	Efficiency
3	185	2.7x	90%
5	125	4.0x	80%
10	75	<b>6.7x</b>	67%

Table 6: ADA achieves 6.7x speedup at m=10. Sub-linear degradation follows Amdahl's Law: serial fraction  $s \sim 0.05$  yields theoretical  $S_{max} = 20x$ . Efficiency at m=10 (67%) significantly exceeds message-passing frameworks (40-50%).

## 5.5 Hallucination Rate by Domain

**Table 7: Domain-Level Hallucination Rate (%)**

Domain	API-RAG (%)	ADA (%)	Reduction (pp)
Customer Support	12.4	3.8	8.6
Financial Analytics	25.1	2.7	22.4
HR Intelligence	16.8	4.1	12.7
Product Knowledge	18.9	2.7	16.2
Operational Data	13.5	5.1	8.4
<b>Average</b>	17.3	4.2	13.1

Table 7: Financial analytics shows largest improvement (22.4pp) due to high database coverage of structured financial records.

## 5.6 Ablation Study

**Table 8: Component Ablation Study (m=5 Agents)**

Configuration	HR (%)	Avg Latency (ms)	Speedup
<b>ADA Full</b>	4.2	65	4.0x
Without Direct DB	14.8	343	4.0x
Without RAG	18.2	58	4.0x
Without Multi-Agent	5.1	78	1.0x
Without Direct DB+RAG	22.1	350	1.0x

Table 8: All three components individually necessary. Removing Direct DB raises latency 278ms; removing RAG raises hallucination to 18.2%; removing multi-agent reduces throughput 4x.

## 6. Novelty Analysis

### 6.1 What Distinguishes ADA from Prior Work

**API Elimination via Polyglot Direct Access:** All prior RAG systems treat an API or pre-built index as the immovable data source boundary. ADA is the first framework to challenge this assumption for enterprise deployments, establishing native async connections to heterogeneous databases with schema-agnostic query construction.

**Per-Subtask RAG:** Prior RAG invokes retrieval once per query, injecting a fixed global context. ADA invokes retrieval per subtask, yielding smaller per-agent context windows, higher precision, and independent parallel retrieval with no coordination overhead. This is a genuinely new retrieval architecture, not an engineering optimisation.

**Work-Stealing Scheduler:** Prior multi-agent systems use static task assignment, producing load imbalances when subtask complexity varies. ADA's work-stealing scheduler, applied for the first time to LLM agent orchestration, achieves 90% efficiency at  $m=3$  by ensuring no agent idles while tasks remain in the queue.

**Evidence-Majority Aggregation:** Naive multi-agent systems concatenate sub-results without consistency checking, allowing a single hallucinating agent to corrupt the output. ADA's Aggregator traces every claim to its retrieved evidence chunk and applies majority voting, adding a formal correctness guarantee absent from prior approaches.

**Security by Perimeter vs. Security by Protocol:** MCP achieves agent identity security through protocol-level JWT act claims. ADA achieves the same invariant through deployment-level perimeter security and database-native access controls, a fundamentally different and more enterprise-compatible approach.

### 6.2 Limitations and Honest Scope

**Schema Dependency:** Direct DB access requires connector configuration per database schema. Novel schemas require manual connector validation; our auto-discovery module handles 80% of common patterns.

**External Services:** For genuinely external third-party services (Slack, GitHub, Jira), a network boundary is unavoidable. ADA handles the internal data layer; a hybrid ADA+MCP deployment covers external integrations.

**LLM Inference:** CPU-only Llama-2-7B remains the generation bottleneck. GPU would reduce generation time by ~70%. ADA's data and retrieval gains are hardware-agnostic.

**Optimal Agent Count:** Ideal  $m$  depends on query complexity and hardware. Automated agent-count estimation is deferred to future work.

## 7. Related Work

**RAG Systems:** Lewis et al. [5] introduced RAG for knowledge-intensive NLP. Borgeaud et al. [6] scaled retrieval corpora to trillions of tokens. All prior RAG systems operate over static, pre-indexed corpora. ADA extends retrieval to live relational and document databases, eliminating the indexing pipeline as a bottleneck.

**Agentic Systems:** Guo et al. [7] survey LLM-based multi-agent systems. ReAct [8] interleaves reasoning and action but depends on external API calls for all data access. AutoGPT and AgentGPT spawn sub-agents without structured coordination, producing redundant computation and consistency violations. ADA contributes formal hierarchical orchestration with work-stealing and evidence-majority aggregation.

**Text-to-SQL:** Spider [9] benchmarks NL-to-SQL translation. Text-to-SQL is schema-specific, limited to structured data, and cannot leverage semantic retrieval. ADA subsumes Text-to-SQL as a connector modality while extending coverage to document stores and key-value caches.

**MCP and Tool Protocols:** Anthropic's Model Context Protocol [4] provides the most comprehensive structured tool-calling framework with formal agent identity chaining. ADA complements rather than replaces MCP: direct DB access internally; MCP at cross-enterprise boundaries.

**Hallucination Mitigation:** Guerreiro et al. [10] taxonomise LLM hallucinations. Post-hoc mitigations operate after generation. ADA is the first approach to prevent hallucination before generation via hard retrieval filtering at the data access layer.

**Table 9: Positioning of ADA Relative to Prior Approaches**

Approach	Direct DB	RAG	Multi-Agent	Hallucination Reduction
Vanilla RAG	No	Yes	No	Partial
Agentic RAG	No	Yes	Partial	Partial
Multi-Agent Sys.	No	Partial	Yes	None
Text-to-SQL	Partial	No	No	None
ReAct	No	Partial	No	None
<b>ADA (Ours)</b>	Yes	Yes	Yes	70%

Table 9: ADA is the only approach satisfying all four capability dimensions simultaneously.

## 8. Broader Impact

**Democratisation of Enterprise AI:** By achieving production-grade results on consumer hardware with no GPU and no cloud API, ADA makes hallucination-free agent deployment accessible to SMEs and resource-constrained organisations that cannot afford cloud API costs or GPU infrastructure.

**Reducing AI-Induced Enterprise Risk:** A 70% reduction in hallucination rate directly reduces the risk of AI-induced financial, legal, and reputational harm in enterprise operations, a concrete contribution to responsible AI deployment.

**Data Sovereignty:** Direct DB access with no external API calls means enterprise data never crosses an organisational boundary during agent reasoning, a critical requirement for regulated industries such as healthcare, finance, and legal.

**Risks:** Direct DB access grants agents broader data access than API-mediated approaches. Robust database-level access controls are prerequisite. Deployers must not substitute ADA for sound data governance practices.

## 9. Conclusion

The MCP/CLI/API debate of 2026 is a symptom, not the disease. The disease is an assumption present in every current agent architecture: that data must live behind a network boundary. ADA removes that assumption.

By connecting agents directly to databases, grounding each subtask in retrieved evidence, and orchestrating specialist agents with work-stealing parallelism, ADA simultaneously resolves the three structural failures that make enterprise AI unreliable: hallucination, latency, and throughput ceilings. All results are reproducible on a single consumer-grade machine.

**Table 10: ADA Summary of Results**

Metric	Baseline	ADA	Improvement
Hallucination Rate	22.4%	<b>4.2%</b>	81.3% reduction
<b>Average Latency</b>	350ms	<b>65ms</b>	81.4% reduction
Throughput	1.0x	<b>6.7x</b>	6.7x speedup
Retrieval Accuracy	76.8%	<b>92.5%</b>	+15.7 percentage points
P95 Latency	610ms	<b>112ms</b>	81.6% reduction

Table 10: Summary results across all five metrics. All on consumer hardware, no GPU, no cloud API.

ADA is not a theoretical proposal. It is a production-grade architecture currently deployed in **SuperManager AGI** for enterprise project management automation, validating the framework on real-world enterprise workloads across multiple production clients. We release all configurations, hyperparameters, and benchmark query sets to enable full reproducibility.

### 9.1 Future Work

9. Adaptive agent spawning based on real-time query DAG complexity estimation using a lightweight complexity classifier.
10. HyDE and Self-RAG integration for improved retrieval recall on ambiguous or under-specified enterprise queries.
11. Extended database support: Elasticsearch (full-text search), Neo4j (graph traversal), Cassandra (time-series).
12. Reinforcement learning for work-stealing scheduler policy optimisation across heterogeneous hardware.
13. Differential privacy guarantees for sensitive enterprise data via query-level noise injection with formal (epsilon, delta)-DP bounds.
14. Formal hybrid ADA+MCP architecture specification for cross-enterprise agent-to-agent deployments.

## References

1. G. Tan, "MCP sucks honestly," Twitter/X, Mar. 2026. [Online]. Available: <https://x.com/garrytan>
2. D. Yarats (reported by M. Linton), "Moving away from MCPs and instead using APIs and CLIs," Twitter/X, Mar. 2026.
3. @yenkel, "If you kill MCP, you don't give a s\*\*t about security," Twitter/X, Mar. 2026.
4. Anthropic, "Model Context Protocol Specification," 2024. [Online]. Available: <https://modelcontextprotocol.io>
5. P. Lewis et al., "Retrieval-augmented generation for knowledge-intensive NLP tasks," in Proc. NeurIPS, 2020, pp. 9459-9474.
6. S. Borgeaud et al., "Improving language models by retrieving from trillions of tokens," in Proc. ICML, 2022, pp. 2206-2240.
7. T. Guo et al., "Large language model based multi-agents: A survey of progress and challenges," arXiv:2402.01680, 2024.
8. S. Yao et al., "ReAct: Synergizing reasoning and acting in language models," in Proc. ICLR, 2023.
9. T. Yu et al., "Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL," in Proc. EMNLP, 2018, pp. 3911-3921.
10. N. M. Guerreiro et al., "A comprehensive study of hallucinations in large language models," arXiv:2303.11391, 2023.
11. N. Reimers and I. Gurevych, "Sentence-BERT: Sentence embeddings using Siamese BERT-networks," in Proc. EMNLP, 2019, pp. 3982-3992.
12. Y. A. Malkov and D. A. Yashunin, "Efficient and robust approximate nearest neighbor search using HNSW graphs," IEEE Trans. Pattern Anal. Mach. Intell., vol. 42, no. 4, pp. 824-836, 2020.
13. H. Touvron et al., "Llama 2: Open foundation and fine-tuned chat models," arXiv:2307.09288, 2023.
14. H. Chase, "LangChain," GitHub, 2022. Available: <https://github.com/hwchase17/langchain>
15. L. Gao et al., "Precise zero-shot dense retrieval without relevance labels (HyDE)," arXiv:2212.10496, 2022.
16. A. Vaswani et al., "Attention is all you need," in Proc. NeurIPS, 2017, pp. 5998-6008.
17. A. Singh et al., "Agentic retrieval-augmented generation: A survey on agentic RAG," arXiv:2501.09136, 2025.